

# A Rollback in the History of Communication-Induced Checkpointing

Islene C. Garcia, Gustavo M. D. Vieira, and Luiz E. Buzato

**Abstract**—The literature on communication-induced checkpointing presents a family of protocols that use logical clocks to control whether forced checkpoints must be taken. For many years, HMNR, also called Fully Informed (FI), was the most complex and efficient protocol of this family. The Lazy-FI protocol applies a lazy strategy that defers the increase of logical clocks, resulting in a protocol with better performance for distributed systems where processes can take basic checkpoints at different, asymmetric, rates. Recently, the Fully Informed aNd Efficient (FINE) protocol was proposed using the same control structures as FI, but with a stronger and, presumably better, checkpoint-inducing condition. FINE and its lazy version, called Lazy-FINE, would now be the most efficient checkpointing protocols based on logical clocks. This paper reviews this family of protocols, proves a theorem on a condition that must be enforced by all stronger versions of FI, and proves that both FINE and Lazy-FINE do not guarantee the absence of useless checkpoints. As a consequence, FI and Lazy-FI can be rolled back to the position of most efficient protocols of this family of index-based checkpointing protocols.

**Index Terms**—Reliability, Checkpointing/restart, Fault-tolerance

## 1 INTRODUCTION

CHECKPOINTING is a widely used technique that provides fault-tolerance to distributed systems. A local checkpoint is a state of a process that can be recovered after a crash. A consistent global checkpoint [1] is a set of local checkpoints that can be used to recover a system after a global failure. If processes take checkpoints at their own pace, a consistent global checkpoint may not be formed, and, in the worst case, the application may need to rollback to its initial state after a failure. This is the well-known *domino effect* [2] caused by the presence of *useless* checkpoints [3].

Some checkpointing protocols avoid useless checkpoints by using a coordinator and control messages [1], [4]. Others use a communication-induced approach: processes can take *basic* checkpoints autonomously, but the protocol uses information obtained from the exchange of messages among the processes to induce *forced* ones and, thus, to eliminate the occurrence of useless checkpoints [5], [6]. Checkpoint-inducing conditions based on information stored in local variables and messages received are used to control whether a forced checkpoint must be taken before delivering the payload of a message to the application. Therefore, communication-induced protocols are often compared in terms of the number of forced checkpoints and the size of the state (data structures) maintained by each process to support the

decision to take a forced checkpoint. The least the number of forced checkpoints taken and the size of the data structures used the better.

Communication-induced index-based checkpointing protocols implement a variant of Lamport's logical clock [7] to state checkpoint-inducing conditions. Protocols that use this approach, such as [8], [9], [10], and [11], enforce an easily observable property that guarantees that checkpoints stamped with the same clock value form a consistent global checkpoint [10]. Furthermore, index-based protocols have presented better performance than protocols based on the tracking of specific checkpoint patterns [12], [13].

Unfortunately, no matter what mechanism is used to trigger forced checkpoints, there is not an optimal checkpointing protocol for all checkpoint and communication patterns [14]. However, for a particular family of protocols, a stronger (more restrictive) condition always produces a protocol that takes fewer checkpoints than a protocol based on a weaker condition [14]. Evidence obtained from experimental comparisons of checkpointing protocols also indicate that stronger conditions usually lead to more efficient protocols [10], [15].

For many years, the HMNR protocol [10] implemented the strongest index-based checkpoint-inducing condition. This protocol has also been called Fully Informed (FI) [11], because it propagates detailed information about the causal past of the processes. Eventually, the literature began to show efforts to produce further optimized versions of FI. The Lazy-FI [16] approach applies the lazy strategy [17] to increment logical clocks of FI. The Fully Informed aNd Efficient (FINE) protocol [15], [18] is based on a checkpoint condition stronger than the one defined for

- I. C. Garcia and L. E. Buzato are with the Institute of Computing, University of Campinas, Brazil.  
E-mail: {islene, buzato}@ic.unicamp.br
- G. M. D. Vieira is with the Department of Computing at Sorocaba, CCGT, Federal University of São Carlos, Brazil.  
E-mail: gdvieira@ufscar.br

Manuscript received [date]; revised [date].

FI but using the same control information maintained by FI. A lazy version of this protocol, called Lazy-FINE was also proposed [19]. The S-FI [20] protocol aims to take the same number of forced checkpoints as FI, but employing a reduced amount of information per message exchanged, an improvement that makes the protocol more scalable. The DCFI [21] delays non-forced checkpoints in order to reduce the total number of checkpoints in the system.

The contributions of the paper are three. Firstly, it reviews the FI [10] and Lazy-FI [16] checkpointing protocols to single out the similarities in logical structure of the conditions used by them to trigger forced checkpoints. Secondly, it proves a theorem that shows that checkpoint-inducing condition of FI cannot be strengthened without respecting the timestamping rules that guarantee the absence of useless checkpoints. Thirdly, it shows that FINE [15] and Lazy-FINE [19] fail the test established by the theorem, thus causing a rollback in the history of communication-induced checkpointing protocols: FI and Lazy-FI are back as the most efficient protocols of this family of index-based checkpointing protocols.

The rest of the paper is structured as follows. Section 2 presents fundamental concepts. Section 3 addresses index-based checkpointing, describing FI [10] and Lazy-FI [16]. Section 4 presents FINE [15] and Lazy-FINE [19], the theorem about the correctness of FI optimizations, and the checkpoint scenarios that show that these protocols may lead to useless checkpoints. Finally, Section 5 concludes the paper.

## 2 FUNDAMENTAL CONCEPTS

This section defines the meaning of distributed computation, checkpoint, consistent global checkpoint and the mechanisms used to track whether checkpoints belong or not to a consistent global checkpoint.

### 2.1 Distributed computation

A set of  $n$  processes ( $P_1, \dots, P_n$ ) that communicate strictly via unicast messages forms a distributed computation. The communication graph is complete, the channels are reliable, but the transmission delays are unpredictable. There is no global clock or shared memory.

Every process  $P_i$  starts with an event  $e_{i,1}$  and executes a possibly infinite sequence of events ( $e_{i,1}, e_{i,2}, \dots$ ). An internal event can only influence the state of the process that has executed it. External events can be the sending or the receiving of messages. Given global time absence, events can be ordered using solely the notion of “cause-and-effect” enabled by the flow of information generated by the occurrence of internal and external events. Thus, causality can be captured by the *causally precedes* relation over events of a distributed computation [7].

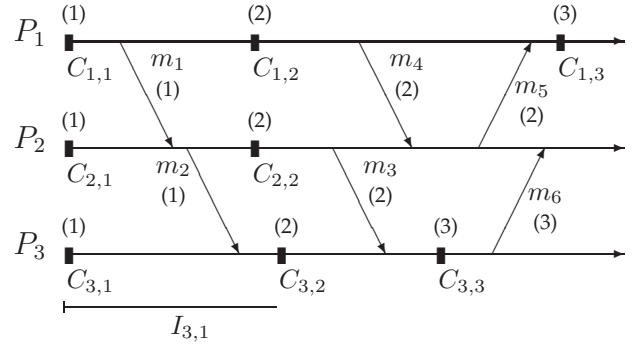


Fig. 1. Processes, checkpoints, and logical clocks

**Definition 2.1** (Causal precedence). *Event  $e_{i,x}$  causally precedes  $e_{j,y}$  ( $e_{i,x} \rightarrow e_{j,y}$ ) if*

- $i = j$  and  $y = x + 1$ , or
- $\exists m : e_{i,x} = \text{send}(m)$  and  $e_{j,y} = \text{receive}(m)$ , or
- $\exists e_{k,z} : e_{i,x} \rightarrow e_{k,z} \wedge e_{k,z} \rightarrow e_{j,y}$ .

### 2.2 Fault tolerance and checkpoints

We assume the crash-recover fault model, that is, in the case of a failure, a process halts and loses its volatile state. During correct execution, processes frequently save their states to stable storage to make possible the recovery of a process (system) by way of a rollback to an earlier consistent state in the case of partial or total system failure.

A checkpoint is the local state of a process that was saved on stable storage. Every process  $P_i$  has an initial checkpoint, denoted by  $C_{i,1}$ , and other checkpoints can be saved along the computation. The  $x$ -th checkpoint of a process  $P_i$  is denoted by  $C_{i,x}$ . An interval  $I_{i,x}$  is the set of events from  $C_{i,x}$  to  $C_{i,x+1}$ , including  $C_{i,x}$  but excluding  $C_{i,x+1}$ .

Let us assume that a process  $P_i$  manages a logical clock  $lc_i$  that is used to timestamp a checkpoint  $C_{i,x}$  with  $C_{i,x}.t$  and a message  $m$  with  $m.t$ . The following rules guarantee that if  $C_{i,x} \rightarrow C_{j,y}$  then  $C_{i,x}.t < C_{j,y}.t$  [10]. These rules can be seen as a specialization of the Lamport’s clock [7] that increments  $lc_i$  only at the occurrence of checkpoints.

- $P_i$  initializes  $lc_i$  at the beginning of the computation;
- $P_i$  increments  $lc_i$  before it saves a checkpoint  $C$  and sets  $C.t = lc_i$ .
- when  $P_i$  sends a message  $m$ , it piggybacks  $lc_i$  on the message (denoted  $m.t$ );
- when  $P_i$  receives a message  $m$ , it sets  $lc_i$  to  $\max(lc_i, m.t)$ .

Fig. 1 depicts a distributed computation. Horizontal lines represent processes, one line per process. Time flows from left to right. Slanted arrows represent messages. Black rectangles are basic checkpoints. The values of the logical clocks—timestamps—associated with each event of interest to the distributed computation are depicted as integers between parentheses.

A checkpoint interval is represented by a left-closed right-open line segment, for example,  $I_{3,1}$ .

### 2.3 Consistency

A consistent global checkpoint is formed by a set of checkpoints that are unrelated by causal precedence [1]. In Fig. 1, for example, the set  $\{C_{1,2}, C_{2,2}, C_{3,2}\}$  is a consistent global checkpoint.

**Definition 2.2** (Consistent global checkpoint). *A global checkpoint  $\{C_{1,x_1}, \dots, C_{n,x_n}\}$  is consistent if*

$$\forall i, j : C_{i,x_i} \not\rightarrow C_{j,x_j}.$$

When two causally unrelated checkpoints cannot be part of the same consistent global checkpoint they must be connected by a sequence of messages called a *zigzag path* [3].

**Definition 2.3** (Zigzag path). *A sequence of messages  $\mu = [m_{l1}, \dots, m_{lq}]$  is a zigzag path from  $C_{i,x}$  to  $C_{j,y}$  if*

- $P_i$  sends  $m_{l1}$  after  $C_{i,x}$ , and
- if  $m_{lz}$ ,  $1 \leq z < q$ , is received by  $P_k$ , then  $m_{l,z+1}$  is sent by  $P_k$  in the same or a later checkpoint interval, and
- $m_{lq}$  is received by  $P_j$  before  $C_{j,y}$ .

The existence of a zigzag path from  $C_{i,x}$  to  $C_{j,y}$  is denoted by  $C_{i,x} \xrightarrow{z} C_{j,y}$ . In Fig. 1,  $[m_1, m_2]$  is a zigzag path from  $C_{1,1}$  to  $C_{3,2}$  such that  $C_{1,1}$  causally precedes  $C_{3,2}$ , being an example of a *causal zigzag path*. The zigzag path  $[m_4, m_3]$  is an example of a *non-causal zigzag path* from  $C_{1,2}$  to  $C_{3,3}$ .

**Definition 2.4** (Z-cycle).  $C_{i,x} \xrightarrow{z} C_{i,x}$

A zigzag path from a checkpoint to itself forms a Z-cycle and makes it impossible for this checkpoint to be part of any consistent global checkpoint. A Z-cycle is the exact condition under which a checkpoint becomes useless [3]. In Fig. 1,  $C_{3,3}$  is useless due to the Z-cycles  $[m_6, m_3]$  and  $[m_6, m_5, m_4, m_3]$ .

#### 2.3.1 Z-consistent timestamping

**Definition 2.5** (Z-consistent timestamping). *A timestamping is consistent with the existence of zigzag paths if*

$$C_{i,x} \xrightarrow{z} C_{j,y} \Rightarrow C_{i,x}.t < C_{j,y}.t$$

A Z-consistent timestamping does not admit a Z-cycle, say  $C \xrightarrow{z} C$ , since the relationship  $C.t < C.t$  is impossible using integers as timestamps [10], [22].

Fig. 1 does not present a Z-consistent timestamping, since  $C_{3,3} \xrightarrow{z} C_{1,3}$  and  $C_{3,3}.t = C_{1,3}.t$ . Fig. 2 shows the same basic checkpoint and communication pattern of Fig. 1 augmented with one forced checkpoint, represented by a black diamond. This extra checkpoint allows the Z-consistent timestamping under the rules presented in Section 2.2. In the next section, forced checkpoints induced by checkpointing protocols will guarantee the enforcement of Z-consistent timestamping.

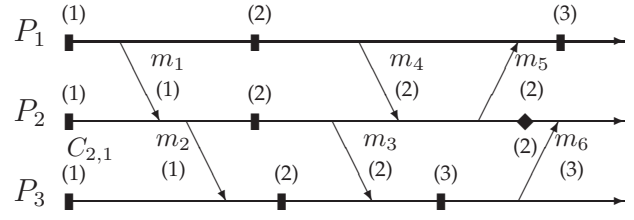


Fig. 2. Z-consistent timestamping

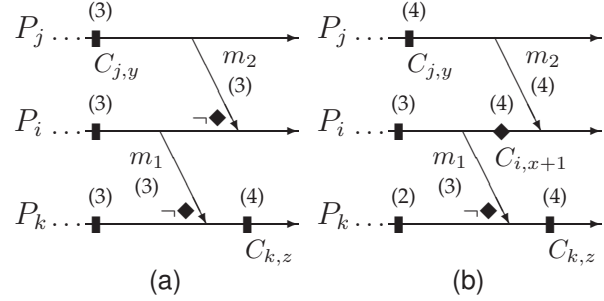


Fig. 3. Partly-informed strategy

## 3 INDEX-BASED CHECKPOINTING

This section starts with the description of a partly-informed strategy to induce forced checkpoints. After that, it presents the fully-informed and lazy strategies describing data structures and checkpoint-inducing conditions that can be used to reduce the number of forced checkpoints in comparison to the partly-informed strategy.

### 3.1 Partly-informed strategy

The partly-informed strategy produces a Z-consistent timestamping by not allowing logical clocks to decrease along a zigzag path. Let us consider zigzag paths composed by two messages  $[m_2, m_1]$ , as the ones depicted in Fig. 3. In Fig. 3a,  $P_i$  receives  $m_2$  and  $m_2.t = m_1.t$ . In Fig. 3b,  $P_i$  receives  $m_2$ , and since  $m_2.t > m_1.t$ ,  $P_i$  takes a forced checkpoint, represented by a black diamond, before delivering  $m_2$  to the application.

This approach can be implemented by the following control structures [10]:

- Boolean array  $sent\_to_i[1 \dots n]$ :  $sent\_to_i[j]$  indicates whether processes  $P_i$  has sent a message to  $P_j$  in the current checkpointing interval.
- Array  $min\_to_i[1 \dots n]$ :  $min\_to_i[j]$  indicates the timestamp of the first message sent in the current interval by  $P_i$  to  $P_j$  or  $+\infty$  if no such message has been sent.

The partly-informed checkpoint-inducing condition  $\mathcal{C}_{PI}$ , evaluated by process  $P_i$  when it receives a message  $m$ , can be expressed as [10]:

$$\mathcal{C}_{PI} \equiv \exists k : sent\_to_i[k] \wedge m.t > min\_to_i[k]$$

As a consequence of the partly-informed strategy,  $P_k$  does not need to take a checkpoint before the reception of  $m_1$ , even when  $m_1.t$  is greater than  $lc_k$ , as in Fig. 3b.

### 3.2 Fully-informed strategy

Taking into account a zigzag path  $[m_2, m_1]$  from  $C_{j,y}$  to  $C_{k,z}$ , like the ones depicted on Fig. 3, the fully-informed strategy [10] explores  $P_i$ 's information about  $C_{k,z}.t$  to establish if Z-consistent timestamping is being preserved. This information is propagated by piggybacking timestamp vectors and checkpoint vectors that carry causal information about  $P_k$  to  $P_i$ .

#### Strengthening the partly-informed strategy

Let us assume that each process  $P_i$  maintains and propagates an array with information about the logical clock of all processes in the computation. Let us define the vector  $clock_i$ , such that  $clock_i[i]$  is equivalent to  $lc_i$  and  $clock_i[k]$  is the highest value of  $lc_k$  that  $P_i$  knows about due to a traditional piggybacking mechanism. Fig. 4 illustrates timestamp vectors, with the logical clocks of processes and messages emphasized in boldface.

In Fig. 4a, there is a message receive where the  $\mathcal{C}_{PI}$  condition is true.  $P_2$  sends a message  $m_1$  to  $P_3$  with  $m_1.t = 1$  and it receives  $m_3$  from  $P_1$  with  $m_3.t = 2$ . However, since  $m_3.clock[3] = 2$  it does not need to save a forced checkpoint. In this case,  $C_{1,2} \xrightarrow{z} C_{3,3}$ ,  $C_{1,2}.t < C_{3,3}.t$  and a Z-consistent timestamping is enforced.

In Fig. 4b, when  $P_2$  receives  $m_2$  from  $P_3$  a similar scenario occurs. The partly-informed condition is true, since  $m_1.t = 1$  and  $m_2.t = 2$ . However,  $m_2.clock[3] = 2$  and  $P_2$  does not need to save a forced checkpoint. When  $P_2$  receives  $m_3$  from  $P_1$  with  $m_3.t = 2$ , the partly-informed condition is true again. Nevertheless, since  $clock_2[3] = 2$  a forced checkpoint is not necessary. As in the previous example,  $C_{1,2} \xrightarrow{z} C_{3,3}$ ,  $C_{1,2}.t < C_{3,3}.t$  and a Z-consistent timestamping is enforced.

A variation  $\mathcal{C}_{FI-1}$  of the partly-informed checkpoint-inducing condition, evaluated by process  $P_i$  when it receives a message  $m$  and that takes into account the value of  $lc_k$  up to  $P_i$ 's knowledge, can be expressed as [10]:

$$\mathcal{C}_{FI-1} \equiv \exists k : sent\_to_i[k] \wedge m.t > min\_to_i[k] \wedge m.t > max(clock_i[k], m.clock[k])$$

The  $\mathcal{C}_{FI-1}$  condition can also be implemented with a reduced set of data structures [10] to minimize the cost of piggybacking information about the processes causal past. In the FI algorithm, there is no need for a process to know exactly which is the clock of another process. It is just important to know if their clocks are synchronized or not. Thus, FI can be rewritten using  $lc$

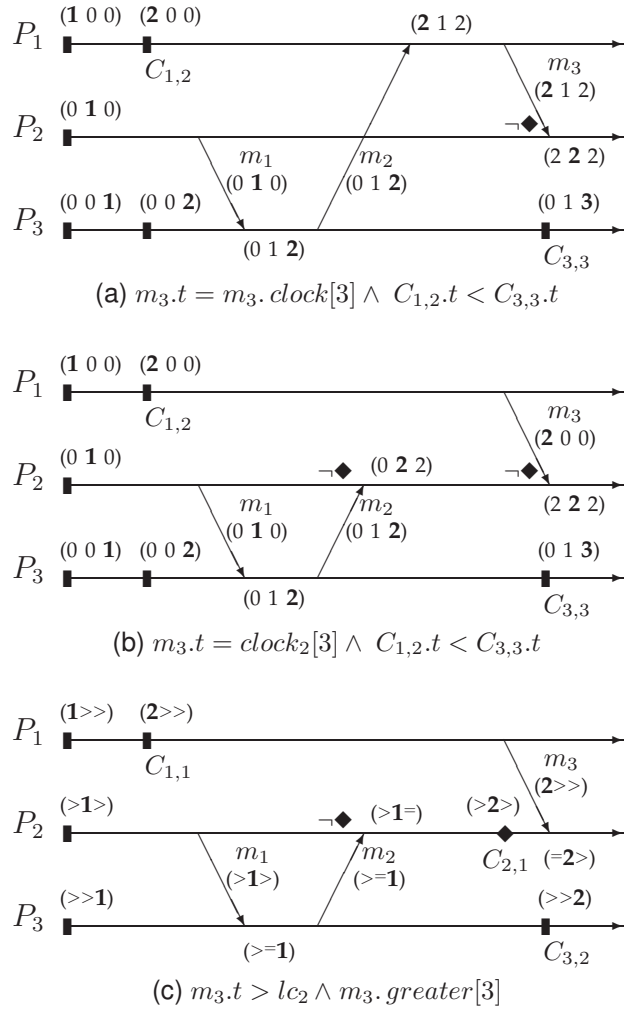


Fig. 4. Strengthening the partly-informed strategy

and a vector of booleans *greater* instead of the vector of integers *clock*.

Each entry  $greater_i[k]$  is true if to the knowledge of  $P_i$  its clock is greater than the clock of  $P_k$  ( $greater_i[k] \equiv clock_i[i] > clock_i[k]$ ). When an entry  $greater_i[k]$  is false,  $clock_i[i]$  must be equal to  $clock_i[k]$ ; their clocks are synchronized. Due to the update rules of the timestamps, it is not possible that  $clock_i[i] < clock_i[k]$ . Fig. 4c presents a distributed computation with *greater* vectors, but instead of using true and false to indicate the truthfulness/falseness of the predicate we have used the signs  $>$  or  $=$ . The logical clock of processes and messages are emphasized in boldface and they are maintained in their positions in the vectors.

When  $P_2$  receives  $m_2$  from  $P_3$  with  $m_2.t = 1$  it trivially does not need to take a forced checkpoint. When  $P_2$  receives  $m_3$  from  $P_1$  with  $m_3.t = 2$ , since  $m_3.greater[3]$ ,  $P_2$  takes a forced checkpoint before delivering  $m_3$ . The FI condition can be stated as [10]:

$$\mathcal{C}_{FI-1} \equiv \exists k : sent\_to_i[k] \wedge m.greater[k] \wedge m.t > lc_i$$



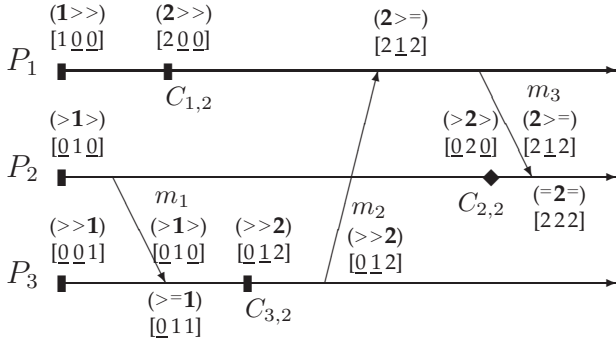


Fig. 5. Fully-informed strategy

### Breaking $[\mu, m]$ Z-cycles

Unfortunately, the  $\mathcal{C}_{FI,1}$  condition strengthens the partly-informed strategy more than it should. In Fig. 5,  $P_2$  receives  $m_3$  from  $P_1$  with  $m_3.t = 2$  and  $m_3.greater[3]$  indicates that  $lc_3$  has already reached 2. However, if  $P_2$  had not taken a checkpoint before delivering  $m_3$ , a Z-cycle  $[m_2, m_3, m_1]$  would have been formed and  $C_{3,2}$  could have become a useless checkpoint. Checkpoint vectors with taken marks can be used to prevent this Z-cycle and any one composed by a causal component  $\mu$  and a single message  $m$ .

Let us assume that each process  $P_i$  maintains and propagates a variation of the traditional vector

clock [23] that counts how many checkpoints have been taken during the computation. The entry  $ckpt_i[i]$  expresses exactly the number of checkpoints taken by  $P_i$  and  $ckpt_i[k]$  counts how many checkpoints  $P_k$  has taken to the best knowledge of  $P_i$ .

An extra boolean array  $taken_i$ <sup>1</sup> can be used to indicate if the causal components ending in the current interval contain a checkpoint. An entry  $taken_i[k]$  is true if there is a causal zigzag path from  $C_{k,ckpt_i[k]}$  to  $C_{i,ckpt_i[i]+1}$  and this causal zigzag path includes a checkpoint [10]. Fig. 5 shows a computation with *greater* vectors and checkpoint vectors with taken marks; an entry  $k$  of the checkpoint array is underlined only if  $taken[k]$  is true.

The condition to break  $[\mu, m]$  Z-cycles can be expressed using the following condition  $\mathcal{C}_{FI,2}$ , evaluated by process  $P_i$  when it receives a message  $m$  [10]:

$$\mathcal{C}_{FI,2} \equiv m.ckpt[i] = ckpt_i[i] \wedge m.taken[i]$$

The checkpoint-inducing condition of FI is an or-operation of the conditions explained before:

$$\mathcal{C}_{FI} \equiv \mathcal{C}_{FI,1} \vee \mathcal{C}_{FI,2}$$

Fig. 6 presents the code that implements FI [10].

1. The *taken* array has the opposite meaning of the boolean array *simple<sub>i</sub>* ( $taken_i[j] \equiv \neg simple_i[j]$ ), defined in the context of a checkpointing protocol [24] that enforces the Rollback Dependency Trackability property [25].

```

take_checkpoint():
  for all k do
    sent_to_i[k] ← false;
  end for
  for all k ≠ i do
    taken_i[k] ← true;
    greater_i[k] ← true;
  end for
  lc_i ← lc_i + 1;
  Save the current state on stable memory;
  ckpt_i[i] ← ckpt_i[i] + 1;

P_i's initialization:
  for all k do
    ckpt_i[k] ← 0;
  end for
  lc_i ← 0;
  taken_i[i] ← false;
  greater_i[i] ← false;
  take_checkpoint();

P_i sends a message to P_k:
  sent_to_i[k] ← true;
  send(m, lc_i, greater_i, ckpt_i, taken_i) to P_k;

FI_1():
  return ∃k : sent_to_i[k] ∧ m.greater[k] ∧ m.t > lc_i;

FI_2():
  return m.ckpt[i] = ckpt_i[i] ∧ m.taken[i];

P_i receives a message from P_j:
  if FI_1() ∨ FI_2() then
    take_checkpoint();
  end if
  if m.t > lc_i then
    lc_i ← m.t
    for all k ≠ i do
      greater_i[k] = m.greater[k];
    end for
  else if m.t = lc_i then
    for all k ≠ i do
      greater_i[k] = greater_i[k] ∧ m.greater[k];
    end for
  end if
  for all k ≠ i do
    if m.ckpt[k] > ckpt_i[k] then
      ckpt_i[k] ← m.ckpt[k];
      taken_i[k] ← m.taken[k];
    else if m.ckpt[k] = ckpt_i[k] then
      taken_i[k] ← taken_i[k] ∨ m.taken[k];
    end if
  end for
  deliver(m)

```

Fig. 6. FI protocol [10]

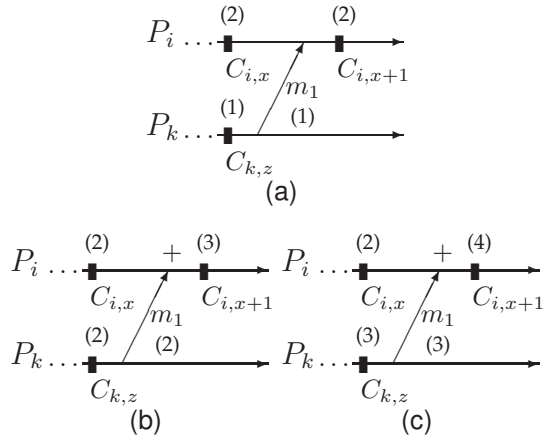


Fig. 7. Lazy indexing strategy

### 3.3 Lazy strategy

The lazy strategy reduces the number of forced checkpoints necessary to produce a Z-consistent timestamping by detecting when a basic checkpoint can be taken by a process  $P_i$  without incrementing its logical clock  $lc_i$ . Let us consider the situation of a process  $P_i$  that receives a single message  $m_1$  in a checkpoint interval and later decides to take a basic checkpoint that ends this interval, as depicted in Fig. 7. If  $m_1.t < C_{i,x}.t$ ,  $P_i$  can reuse the same timestamp of  $C_{i,x}$  to label  $C_{i,x+1}$  because  $C_{k,z}.t < C_{i,x+1}.t$  will still hold, as Fig. 7a shows. However, if  $m_1.t \geq C_{i,x}.t$ ,  $P_i$  must increment  $lc_i$  to label  $C_{i,x+1}$  in order to produce a Z-consistent timestamping where  $C_{k,z}.t < C_{i,x+1}.t$ , as depicted in Fig. 7b and Fig. 7c.

The lazy strategy can be implemented by introducing a flag *increment* that signals  $P_i$  it must increment  $lc_i$  before taking a basic checkpoint. This flag is set to false in the beginning of each checkpoint interval and is set to true whenever a message with  $m.t \geq lc_i$  is received. The setting of the *increment* flag is illustrated in Fig. 7 by a + sign.

#### Lazy-FI protocol

Let us try to apply the lazy approach to FI using the vector *greater*. In Fig. 8a, when  $P_2$  receives  $m_5$  from  $P_1$  with a greater clock, it can verify that  $P_3$  have already reached the same clock. However, due to the lazy strategy,  $P_2$  does not know whether  $P_3$  will increase its clock to save the next checkpoint. Thus, a forced checkpoint before the delivering of  $m_5$  will be required in order to guarantee a Z-consistent timestamping.

Fortunately, a small variation in the *greater* vector allows the implementation of the lazy strategy. The entry  $greater_i[i]$  is set to false at every checkpoint and set to true when a process set its *increment* flag [16]. To differentiate the vector used in Lazy-FI to the one used in FI, we are going to introduce an equivalent vector with an intuitive meaning: *equal\_incr*. Each entry  $equal\_incr_i[k]$  is true if to the knowledge of

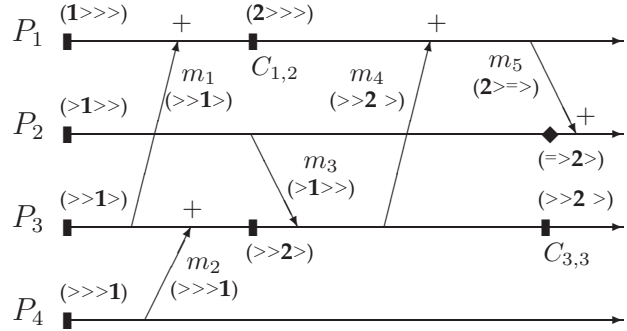
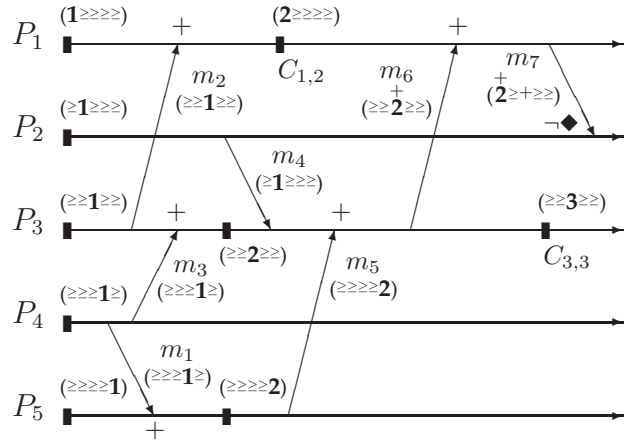
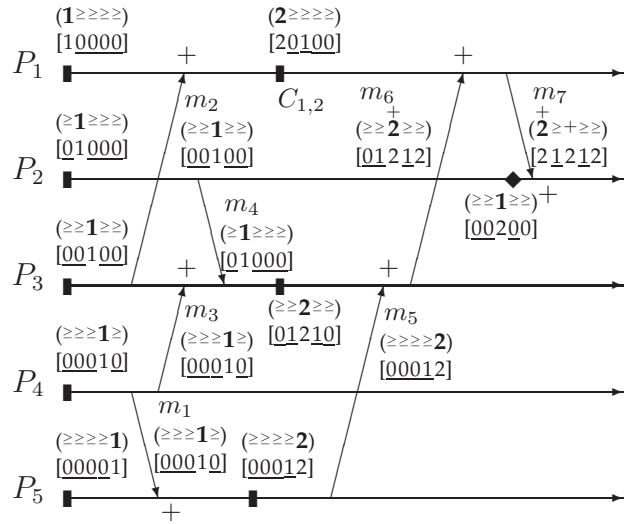
(a) *lc* and *greater* are not enough(b) *lc* and *equal\_incr*(c) *lc*, *equal\_incr*, *ckptv*, and *taken*

Fig. 8. Timestamp information to implement Lazy-FI

```

take_checkpoint():
  for all k do
    sent_toi[k] ← false;
  end for
  for all k ≠ i do
    takeni[k] ← true;
  end for
  if increment then
    lci ← lci + 1;
    for all k do
      equal_incri[k] ← false;
    end for
  end if
  increment ← false;
  Save the current state on stable memory;
  ckpti[i] ← ckpti[i] + 1;

Pi's initialization:
  for all k do
    ckpti[k] ← 0;
  end for
  lci ← 0;
  increment ← true;
  takeni[i] ← false;
  take_checkpoint();

Pi sends a message to Pk:
  sent_toi[k] ← true;
  send(m, lci,
    equal_incri, ckpti, takeni) to Pk;

LAZY_FI_1():
  return ∃k : sent_toi[k] ∧ ¬m.equal_incr[k] ∧ m.t > lci;

LAZY_FI_2():
  return m.ckpt[i] = ckpti[i] ∧ m.taken[i];

Pi receives a message from Pj:
  if LAZY_FI_1() ∨ LAZY_FI_2() then
    take_checkpoint();
  end if
  if m.t > lci then
    lci ← m.t;
    incrementi ← true;
    equal_incr[i] ← true;
    for all k ≠ i do
      equal_incri[k] ← m.equal_incri[k];
    end for
  else if m.t = lci then
    incrementi ← true;
    equal_incr[i] ← true;
    for all k do
      equal_incri[k] ← equal_incri[k] ∨ m.equal_incri[k];
    end for
  end if
  for all k ≠ i do
    if m.ckpt[k] > ckpti[k] then
      ckpti[k] ← m.ckpt[k];
      takeni[k] ← m.taken[k];
    else if m.ckpt[k] = ckpti[k] then
      takeni[k] ← takeni[k] ∨ m.taken[k];
    end if
  end for
  deliver(m);

```

Fig. 9. Lazy-FI protocol (adapted from [16])

$P_i$  its clock is equal to the clock of  $P_k$  and  $P_k$  will increase its clock before saving the next checkpoint. When  $equal\_incr_i[k]$  is false, we have no additional information whether the clock of  $P_i$  is greater or equal to the clock of  $P_k$ .

Fig. 8b shows the propagation of  $equal\_incr$  and it is very similar to Fig. 8a. Once again, instead of true and false values, we have used the signs + and ≥. Due to an extra message from  $P_5$ ,  $P_3$  will increase its clock before saving the next checkpoint.  $P_2$  receives this information and does not take a forced checkpoint before delivering  $m_7$ .

Fig. 8c shows another similar situation that emphasizes the need of the vectors  $ckpt$  and  $taken$ . Although upon the reception of  $m_7$   $P_2$  receives the information that  $P_3$  will increase its clock,  $P_2$  will take a forced checkpoint to break the Z-cycle  $[m_7, m_4, m_6]$ .

The conditions used in the Lazy-FI protocol can be stated as follows:

$$\mathcal{C}_{Lazy-FI} \equiv \mathcal{C}_{Lazy-FI_1} \vee \mathcal{C}_{Lazy-FI_2}$$

$$\mathcal{C}_{Lazy-FI_1} \equiv \exists k : sent\_to_i[k] \wedge \neg m.equal\_incr[k] \wedge$$

$$m.t > lc_i$$

$$\mathcal{C}_{Lazy-FI_2} \equiv \mathcal{C}_{FI_2}$$

Fig. 9 presents the code that implements the Lazy-FI protocol [16] using the  $equal\_incr$  vector.

## 4 ATTEMPTS TO OPTIMIZE FI AND LAZY-FI

This section starts with a description of the FINE approach to optimize FI. After that, it presents a property that must be followed by all optimizations of FI and proves that both FINE and Lazy-FINE do not guarantee the absence of useless checkpoints.

### 4.1 The FINE approach

The basic FINE protocol tries to reduce the number of forced checkpoints using the same data structures as the FI protocol. Fig. 10 illustrates the approach.  $P_2$  has sent a message  $m_1$  to  $P_3$  with  $m_1.t = 1$ . When  $P_2$  receives  $m_3$  from  $P_1$ , it verifies that  $m_3.t > m_1.t$  and  $lc_3$ , up to  $P_2$ 's knowledge, has not reached 2

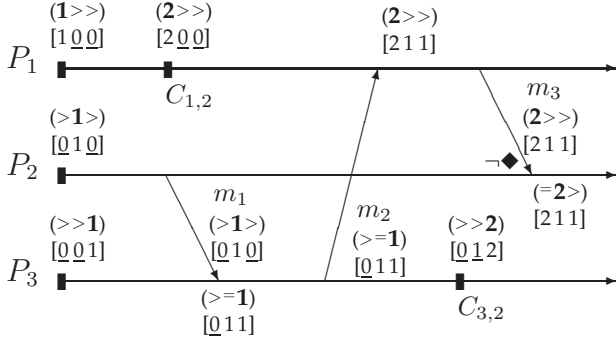


Fig. 10. Fine proposal

**FINE\_10**

```

return  $\exists k : \text{sent\_to}_i[k] \wedge m.\text{greater}[k] \wedge$ 
 $m.t > lc_i \wedge m.\text{taken}[k]$ 

```

Fig. 11. Checkpoint-inducing condition for  $\mathcal{C}_{\text{FINE}_1}$  [19]**Lazy-FINE\_10**

```

return  $\exists k : \text{sent\_to}_i[k] \wedge \neg m.\text{equal\_incr}[k] \wedge$ 
 $m.t > lc_i \wedge m.\text{taken}[k]$ 

```

Fig. 12. Checkpoint-inducing condition for  $\mathcal{C}_{\text{Lazy-FINE}_1}$  (adapted from [19])

yet. However, since  $m_3.\text{taken}[2]$  is false, the messages involve no Z-cycle. They are called harmless cycles.

The basic FINE protocol is based on the following condition [15]:

$$\mathcal{C}_{\text{FINE}} \equiv \mathcal{C}_{\text{FINE}_1} \vee \mathcal{C}_{\text{FINE}_2}$$

where condition  $\mathcal{C}_{\text{FINE}_1}$  can be expressed using a *greater* vector [19] and  $\mathcal{C}_{\text{FINE}_2}$  is equivalent to  $\mathcal{C}_{\text{FI}_2}$ .

$$\mathcal{C}_{\text{FINE}_1} \equiv \exists k : \text{sent\_to}_i[k] \wedge m.\text{greater}[k] \wedge m.t > lc_i \wedge m.\text{taken}[k]$$

$$\mathcal{C}_{\text{FINE}_2} \equiv \mathcal{C}_{\text{FI}_2}$$

Fig. 11 presents the code that implements the  $\mathcal{C}_{\text{FINE}_1}$  predicate. The complete basic FINE protocol is implemented by replacing the **FI\_10** with **FINE\_10** in Fig. 6.

A lazy version of the FINE protocol, called Lazy-FINE, has been proposed in the literature [19]. Let us define the checkpoint inducing conditions using the vector *equal\_incr*:

$$\mathcal{C}_{\text{Lazy-FINE}} \equiv \mathcal{C}_{\text{Lazy-FINE}_1} \vee \mathcal{C}_{\text{Lazy-FINE}_2}$$

$$\mathcal{C}_{\text{Lazy-FINE}_1} \equiv \exists k : \text{sent\_to}_i[k] \wedge \neg m.\text{equal\_incr}[k] \wedge m.t > lc_i[i] \wedge m.\text{taken}[i]$$

$$\mathcal{C}_{\text{Lazy-FINE}_2} \equiv \mathcal{C}_{\text{FI}_2}$$

Fig. 12 presents the code that implements the  $\mathcal{C}_{\text{Lazy-FINE}_1}$  predicate. The complete basic Lazy-FINE protocol is implemented by replacing the **Lazy\_FI\_10** with **Lazy\_FINE\_10** in Fig. 9.

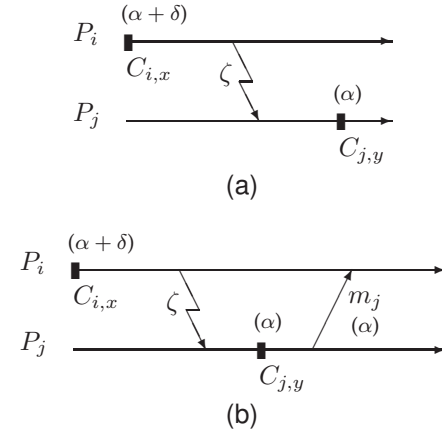


Fig. 13. Necessity of a Z-consistent timestamping

## 4.2 FI's optimization limits

The timestamps of Fig. 10 are not Z-consistent, since  $C_{1,2} \xrightarrow{z} C_{3,2}$  and  $C_{1,2}.t = C_{3,2}.t$ . This violation of Z-consistency may appear innocuous at first, but it violates an important property of any FI optimization. Suppose one considers the  $\mathcal{C}_{\text{FI}_1} \wedge \mathcal{P} \vee \mathcal{C}_{\text{FI}_2}$  as capable of producing a more efficient protocol, and  $\mathcal{C}_{\text{FI}_2}$  is kept exactly as in the original condition. Theorem 4.1 proves that  $\mathcal{C}_{\text{FI}_1} \wedge \mathcal{P}$  must enforce a Z-consistent timestamping to be a valid optimization.

**Theorem 4.1.** *Any optimization of the FI protocol whose checkpoint-inducing condition can be expressed as  $\mathcal{C}_{\text{FI}_1} \wedge \mathcal{P} \vee \mathcal{C}_{\text{FI}_2}$ ,  $\mathcal{C}_{\text{FI}_1} \wedge \mathcal{P}$  must enforce a Z-consistent timestamping.*

*Proof:* Assume an optimization of the FI protocol with condition  $\mathcal{C}_{\text{FI}_1} \wedge \mathcal{P} \vee \mathcal{C}_{\text{FI}_2}$  where  $\mathcal{C}_{\text{FI}_1} \wedge \mathcal{P}$  does not enforce a Z-consistent timestamping. Thus, there must be a computation with two checkpoints  $C_{i,x}$  and  $C_{j,y}$  such as  $C_{i,x} \xrightarrow{z} C_{j,y}$  and  $C_{i,x}.t \geq C_{j,y}.t$ . For simplicity, let  $C_{i,x}.t = \alpha + \delta$ ,  $C_{j,y}.t = \alpha$  with  $\delta \geq 0$  and  $\zeta$  be the zigzag path between  $C_{i,x}$  and  $C_{j,y}$  (Fig. 13a). Depending on the properties of  $\mathcal{C}_{\text{FI}_1} \wedge \mathcal{P}$ , this computation can be arbitrarily complex, involving other processes and requiring many messages exchanges. Let us assume that this computation does not enforce a Z-consistent timestamping, but has no useless checkpoint.

We add to the computation another message  $m_j$  sent by  $P_j$  as the *first* event of the interval  $I_{j,y}$  and received by  $P_i$  in the interval  $I_{i,x}$  after the first message of  $\zeta$  is sent (Fig. 13b). By our construction  $m_j.t = \alpha$ . This implies that  $C_{i,x} \nrightarrow C_{j,y}$ , and that  $\mathcal{C}_{\text{FI}_2}$  can never be true. Evaluating  $\mathcal{C}_{\text{FI}_1}$  at the time  $m_j$  is received,  $m_j.t \leq \alpha + \delta$  and no checkpoint is forced upon the reception of  $m_j$ , no matter the existence of  $\mathcal{P}$ . Since  $\zeta$  must be non-causal, the z-cycle formed by  $m_j$  and  $\zeta$  is not detected by  $\mathcal{C}_{\text{FI}_2}$  and the protocol allowed the occurrence of a useless checkpoint  $C_{j,y}$ .  $\square$



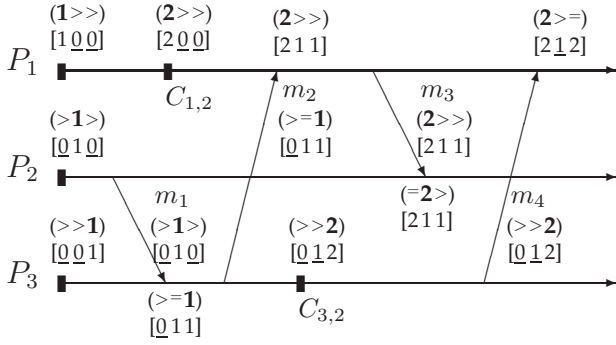


Fig. 14. A useless checkpoint under FINE

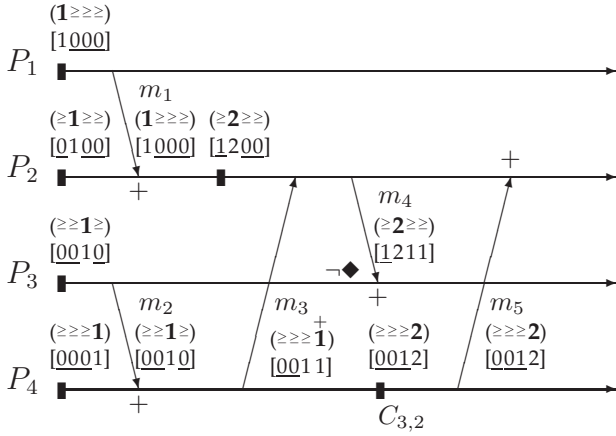


Fig. 15. A useless checkpoint under Lazy-FINE

#### *FINE may lead to useless checkpoints*

According to Theorem 4.1, the  $\mathcal{C}_{FINE\_1}$  predicate produces a checkpointing protocol that does not guarantee the absence of useless checkpoints. Indeed, Fig. 14 is a counterexample: it shows a possible continuation of the scenario presented in Fig. 10 that leads to the occurrence of a useless checkpoint. When  $P_2$  receives  $m_3$  from  $P_1$ , there is no Z-cycle known to  $P_2$  closed by the receipt of  $m_3$ . However, this does not exclude the formation of a Z-cycle, undetected at the time  $m_3$  is received.

#### *Lazy-FINE may lead to useless checkpoints*

Theorem 4.1 also informs us that the  $\mathcal{C}_{Lazy-FINE\_1}$  predicate does not guarantee the absence of useless checkpoints. In Fig. 15, when  $P_3$  receives  $m_4$  from  $P_2$ ,  $m_4.t > lc_3$ , but since  $\neg m.taken[3]$ ,  $[m_4, m_2, m_3]$  would form just a harmless Z-cycle. However,  $P_3$  receives  $m_5$  from  $P_4$  with  $m_5.t = lc_3$  no forced checkpoint is taken and a Z-cycle  $m_5, m_4, m_2$  is formed.

## 5 CONCLUSION

This paper reviewed index-based checkpointing protocols and proved that the FINE and Lazy-FINE protocols do not guarantee the absence of useless checkpoints. This paper also reinforces that all optimizations of FI must enforce a Z-consistent timestamping.

As a consequence, FI and Lazy-FI can be rolled back to the position of most efficient index-based protocols; whether or not they can be further optimized remains an open problem.

## REFERENCES

- [1] K. M. Chandy and L. Lamport, "Distributed snapshots: determining global states of distributed systems," *ACM Trans. Comput. Syst.*, vol. 3, no. 1, pp. 63–75, Feb. 1985.
- [2] B. Randell, "System structure for software fault tolerance," in *Proceedings of the international conference on Reliable software*, vol. 10, no. 6. New York, NY, USA: ACM, 1975, pp. 437–449.
- [3] R. H. B. Netzer and J. Xu, "Necessary and Sufficient Conditions for Consistent Global Snapshots," *IEEE Transactions on Parallel and Distributed Systems*, vol. 6, no. 2, pp. 165–169, Feb. 1995.
- [4] R. Koo and S. Toueg, "Checkpointing and Rollback-Recovery for Distributed Systems," *IEEE Transactions on Software Engineering*, vol. SE-13, no. 1, pp. 23–31, Jan. 1987.
- [5] E. N. M. Elnozahy, L. Alvisi, Y. M. Wang, and D. B. Johnson, "A survey of rollback-recovery protocols in message-passing systems," *ACM Comput. Surv.*, vol. 34, no. 3, pp. 375–408, Sep. 2002.
- [6] D. Manivannan and M. Singhal, "Quasi-synchronous checkpointing: Models, characterization, and classification," *IEEE Transactions on Parallel and Distributed Systems*, vol. 10, no. 7, pp. 703–713, 1999.
- [7] L. Lamport, "Time, clocks, and the ordering of events in a distributed system," *Commun. ACM*, vol. 21, no. 7, pp. 558–565, Jul. 1978.
- [8] D. Briatico, A. Ciuffoletti, and L. Simoncini, "A Distributed Domino-Effect Free Recovery Algorithm," in *4th IEEE Symposium on Reliability in Distributed Software and Database Systems*, 1984.
- [9] D. Manivannan and M. Singhal, "A low-overhead recovery technique using quasi-synchronous checkpointing," in *Proceedings of 16th International Conference on Distributed Computing Systems*. IEEE Comput. Soc. Press, 1996, pp. 100–107.
- [10] J. M. Hélary, A. Mostefaoui, R. H. B. Netzer, and M. Raynal, "Communication-based prevention of useless checkpoints in distributed computations," *Distributed Computing*, vol. 13, no. 1, pp. 29–43, Jan. 2000.
- [11] J. Tsai, "An Efficient Index-Based Checkpointing Protocol with Constant-Size Control Information on Messages," *IEEE Trans. Dependable Secur. Comput.*, vol. 2, no. 4, pp. 287–296, 2005.
- [12] L. Alvisi, E. Elnozahy, S. Rao, S. A. Husain, and A. de Mel, "An analysis of communication induced checkpointing," in *Fault-Tolerant Computing, 1999. Digest of Papers. Twenty-Ninth Annual International Symposium on*. IEEE, 1999, pp. 242–249.
- [13] G. M. D. Vieira and L. E. Buzato, "Distributed checkpointing: Analysis and benchmarks," in *SBRC '06: Proc. of the 24th Brazilian Symposium on Computer Networks*, Curitiba, Paraná, Brazil, May 2006.
- [14] J. Tsai, Y.-M. Wang, and S.-Y. Kuo, "Evaluations of domino-free communication-induced checkpointing protocols," *Information Processing Letters*, vol. 69, no. 1, pp. 31–37, Jan. 1999.
- [15] Y. Luo and D. Manivannan, "FINE: A Fully Informed aNd Efficient communication-induced checkpointing protocol for distributed systems," *Journal of Parallel and Distributed Computing*, vol. 69, no. 2, pp. 153–167, Feb. 2009.
- [16] J. Tsai, "Applying the Fully-Informed Checkpointing Protocol to the Lazy Indexing Strategy," *Journal of Information Science and Engineering*, vol. 23, pp. 1611–1621, 2007. [Online]. Available: [http://www.iis.sinica.edu.tw/page/jise/2007/200709\\_19.html](http://www.iis.sinica.edu.tw/page/jise/2007/200709_19.html)
- [17] G. M. D. Vieira, I. C. Garcia, and L. E. Buzato, "Systematic Analysis of Index-Based Checkpointing Algorithms using Simulation," in *SCTF '01: Proc. of the IX Brazilian Symposium on Fault-Tolerant Computing*, Florianópolis, Santa Catarina, Brazil, Jan. 2001, pp. 31–42. [Online]. Available: <http://www.ic.unicamp.br/~gdrvieira/publications/>
- [18] Y. Luo and D. Manivannan, "FINE: A Fully Informed aNd Efficient Communication-Induced Checkpointing Protocol," in *Third International Conference on Systems (icons 2008)*. IEEE, Apr. 2008, pp. 16–22.

- [19] —, “Theoretical and experimental evaluation of communication-induced checkpointing protocols in  $F_E$  and  $F_{LAZY-E}$  families,” *Performance Evaluation*, vol. 68, no. 5, pp. 429–445, May 2011.
- [20] A. C. Simon, S. E. P. Hernandez, J. R. P. Cruz, P. Gomez-Gil, and K. Drira, “A scalable communication-induced checkpointing algorithm for distributed systems,” *IEICE TRANSACTIONS on Information and Systems*, vol. 96, no. 4, pp. 886–896, 2013.
- [21] C. Simon, A. Calixto, S. E. P. Hernandez, and J. R. Perez Cruz, “A delayed checkpoint approach for communication-induced checkpointing in autonomic computing,” in *Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE)*, 2013 IEEE 22nd International Workshop on. IEEE, 2013, pp. 56–61.
- [22] J.-M. Hélary, A. Mostéfaoui, and M. Raynal, “Virtual precedence in asynchronous systems: Concept and applications,” in *Distributed Algorithms*, ser. Lecture Notes in Computer Science, M. Mavronicolas and P. Tsigas, Eds. Berlin/Heidelberg: Springer Berlin / Heidelberg, 1997, vol. 1320, ch. 14, pp. 170–184.
- [23] C. Fidge, “Logical Time in Distributed Computing Systems,” *Computer*, vol. 24, no. 8, pp. 28–33, Aug. 1991.
- [24] R. Baldoni, J.-M. Hélary, and M. Raynal, “Rollback-Dependency Trackability: A Minimal Characterization and Its Protocol,” *Information and Computation*, vol. 165, no. 2, pp. 144–173, Mar. 2001.
- [25] Y.-M. Wang, “Consistent global checkpoints that contain a given set of local checkpoints,” *IEEE Transactions on Computers*, vol. 46, no. 4, pp. 456–468, Apr. 1997.